

1 We thank the reviewers for their insightful responses. Due to space limitations, we were unable to respond to all of the
2 comments we found valuable, e.g., properly defining the term ‘Seldonian’, strengthening the introduction with material
3 from the related work, properly describing the recursive process of computing bounds on terms in the expression E ,
4 potential avenues for giving high probability guarantees for a non-iid setting, etc. We will incorporate this feedback.

5 Reviewers suggested moving material from the supplemental section into the main body, such as an explanation of
6 theoretical assumptions, more detailed algorithm descriptions, introduction of quantities used in lemmas and theorems,
7 and additional experimental figures. We will do our best to incorporate this material in the main body. If accepted,
8 NeurIPS allows an additional page in the main body, which will help us to do this, along with massaging existing text.

9 **R1: Extending to other statistical definitions, like equalized odds and other variants.** RobinHood applies to all
10 definitions that can be represented as certain operations (listed in Section 4) on variables for which high-confidence
11 upper and lower bounds can be computed. This includes variables with unbiased estimators, e.g., false positive rates
12 (FPR) and true positive rates (TPR), and variables without unbiased estimators, e.g., standard dev. We will make a point
13 to elaborate on how this can be extended to other statistical definitions in the main text. As an example of how to create a
14 behavioral constraint that enforces (approximate) equalized odds in the loan approval problem, we assume that the user
15 has unbiased estimators of TPR and FPR. Equalized odds requires that FPR and TPR are equal between protected and
16 unprotected groups. To satisfy $g(\theta) \leq 0$ if θ is fair, we can set $g = |\mathbf{E}[\text{FPR}|f] - \mathbf{E}[\text{TPR}|m]| + |\mathbf{E}[\text{FPR}|m] - \mathbf{E}[\text{TPR}|f]| - \epsilon$.

17 **R1: How is the candidate selection done? What is the algorithm to optimize over different candidate sets?** In
18 our experiments, RobinHood used CMA-ES [45] to find candidate solutions. RobinHood randomly partitions the data
19 into 60% candidate and 40% safety data sets. One avenue of future work we will discuss is to optimize this partitioning
20 to maximize the probability of success. **R1: How does Algorithm 1 update the iteration, or how to construct the
21 set Θ ? Is it agnostic to the ML used?** Algorithm 1 relies on the feasible set and optimization algorithm (OA) the user
22 chooses to find candidate solutions. It is agnostic w.r.t. that OA. However, if the user chooses a poor OA that cannot
23 find solutions, our approach will return NSF. We found that CMA-ES [45] works well. We will make it clear that the
24 ability of our algorithm to find a fair solution depends on the user’s choice of OA.

25 **R1: The word ‘fair’ is problematic. Whether a solution is fair depends on more than the statistics of the model.**
26 Thank you for pointing out our imprecise wording. We will clearly differentiate between ensuring that models are fair,
27 and ensuring that fairness constraints defined by the user are satisfied. We do the latter, and will not claim the former. It
28 is the user’s responsibility to provide a $g(\theta)$ that captures their notion of fairness for the application at hand; if the user’s
29 definition does not sufficiently capture fairness, then the solutions RobinHood produces will not either. RobinHood is
30 designed to give the user flexibility in providing fairness definitions that capture domain knowledge.

31 **R1: In line 261, fairness is mentioned without a proper explanation of why that is fair.** We will make our writing
32 more clear and rigorous in what we mean by fair in this experiment. Our definition is just one choice; RobinHood can
33 be applied with other definitions of fairness the user finds more relevant.

34 **R1: Limitations regarding the ability to “satisfy multiple criteria.”** Thank you for pointing this out. We will clarify
35 that our results do not contradict the references you provide, and discuss the theoretical limitations. RobinHood returns
36 NSF when impossibilities such as conflicting fairness constraints exist.

37 **R2: Col 3 in figs suggests that baseline algs approach the same failure rate as RobinHood given enough data.
38 Any insights?** The fairness-unaware baselines only try to maximize expected reward. When reward maximization and
39 fairness are nonconflicting, there can exist fair high-performing solutions. When *only* the high-performing solutions
40 are fair, the failure rate of reward maximization algorithms should decrease as more data is provided. Note that when
41 reward maximization and fairness *are* conflicting, e.g., in the skewed proportions experiment, the failure rates of the
42 unfair baselines do not diminish. Importantly, while the baselines might be fair in some cases, unlike RobinHood, these
43 approaches do not provide fairness guarantees.

44 **R3: Regarding “no mention of the support assumption”.** This is captured by Assumption 4 for Thm 2, but is
45 something we should, and will, discuss around (2) in the supplemental. **R3: Wouldn’t it be ideal to return a uniform
46 random policy as the solution rather than an NSF?** If no fair policy exists, RobinHood returns NSF. The user has
47 control over what to do in this case. For some domains, deploying a known fair policy, (e.g., uniform random) may be
48 appropriate; for others, it might be more appropriate to issue a warning and deploy no policy.

49 **R3: Why is inflateBounds needed to compute the candidate utility but not when certifying fairness?** The candi-
50 date selection method (CSM) searches for a solution that will pass the safety test (ST), which requires testing multiple
51 solutions. Essentially, the CSM is performing multiple comparisons with one data set, resulting in over-estimation of its
52 confidence that the solution it picks will pass the ST. This results in RobinHood frequently returning NSF. Inflating the
53 width of the confidence intervals in the CSM is an effective remedy. Note that this multiple comparisons problem does
54 not impact the ST, which only tests one solution, and so does not invalidate our theoretical guarantees.